

PAC-Bayesian Generalization Bound on Confusion Matrix for Multi-Class Classification*

Emilie Morvant Sokol Koço Liva Ralaivola
 Aix-Marseille Univ., LIF-QARMA, CNRS, UMR 7279, F-13013, Marseille, France
 {firstname.name}@lif.univ-mrs.fr

March 23, 2012

Abstract

In this work, we propose a PAC-Bayes bound for the generalization risk of the Gibbs classifier in the multi-class classification framework. The novelty of our work is the critical use of the *confusion matrix* of a classifier as an error measure; this puts our contribution in the line of work aiming at dealing with performance measure that are richer than mere scalar criterion such as the misclassification rate. Thanks to very recent and beautiful results on matrix concentration inequalities, we derive two bounds showing that the true confusion risk of the Gibbs classifier is upper-bounded by its empirical risk plus a term depending on the number of training examples in each class. To the best of our knowledge, this is the first PAC-Bayes bounds based on confusion matrices.

Keywords: Machine Learning, PAC-Bayes generalization bounds, Confusion Matrix, Concentration Inequality, Multi-Class Classification

1 Introduction

The PAC-Bayesian framework, first introduced in McAllester (1999b), is an important field of research in learning theory. It borrows ideas from the philosophy of Bayesian inference and mix them with techniques used in statistical approaches of learning. Given a family of classifiers \mathcal{F} , the ingredients of a PAC-Bayesian bound are a *prior distribution* \mathfrak{P} over \mathcal{F} , a learning sample S and a *posterior distribution* \mathfrak{Q} over \mathcal{F} . Distribution \mathfrak{P} conveys some prior belief on what are the best classifiers from \mathcal{F} (prior any access to S); the classifiers expected to be the most performant for the classification task at hand therefore have the largest weights under \mathfrak{P} . The posterior distribution \mathfrak{Q} is learned/adjusted using the information provided by the training set S . The essence of PAC-Bayesian results is to bound the risk of the *stochastic* Gibbs classifier associated with \mathfrak{Q} Catoni (2004) —in order to predict the label of an example \mathbf{x} , this predictor first draws a classifier f from \mathcal{F} according to \mathfrak{Q} and then returns $f(\mathbf{x})$.

When specialized to appropriate function space \mathcal{F} and relevant families of prior and posterior distributions, PAC-Bayes bounds can be used to characterize the error of different existing classification methods. An example deals with the risk of methods based upon the idea of the majority vote. We may notice that if \mathfrak{Q} is the posterior distribution, the error of the \mathfrak{Q} -weighted majority vote classifier (which makes a prediction for \mathbf{x} according to $\sum_f f(\mathbf{x})\mathfrak{Q}(f)$) is bounded by twice the error of the Gibbs classifier. If the classifiers from \mathcal{F} the \mathfrak{Q} puts a lot of weight on are good enough, the bound on the risk of the Gibbs classifier can therefore be an informative bound for the \mathfrak{Q} -weighted majority vote. Langford and Shawe-taylor (2002) give a PAC-Bayes bound for Support Vector Machine (SVM), which depends on the margin of the examples. In their study, both the prior and posterior distribution are normal distributions, with different means and variances. Empirical results show that this bound is a good estimator of the risk of SVMs Langford (2005).

PAC-Bayes bounds can also be used to derive new supervised learning algorithms. For example, Lacasse et al. (2007) have introduced an elegant bound on the risk of the majority vote, which holds for any space \mathcal{F} . This bound is used to derive an algorithm, namely MinCq, which achieves empirical results on par with state-of-the-art methods.

*This work was supported in part by the french projects VideoSense ANR-09-CORD-026 and DECODA ANR-09-CORD-005-01 of the ANR in part by the IST Programme of the European Community, under the PASCAL2 Network of Excellence, IST-2007-216886. This publication only reflects the authors' views.

Some other important results are given in Catoni (2007), Seeger and Seeger (2002), McAllester (1999a) and Langford et al. (2001).

In the present paper, we address the multi-class classification problem. Some related works are therefore the multi-class formulations for the Support Vector Machines, such as the frameworks presented in Weston and Watkins (1998), Lee et al. (2004) and Crammer and Singer (2002). As majority vote methods, we can also cite multi-class adaptations of the boosting method called AdaBoost Freund and Schapire (1996), such as the framework given in Mukherjee and Schapire (2011), the AdaBoost.MH/AdaBoost.MR algorithms Schapire and Singer (1999) and the SAMME algorithm Zhu et al. (2009).

The originality of our work is that we consider the *confusion matrix* of the Gibbs classifier as an error measure. We believe that in the multi-class framework, it is more relevant to consider the confusion matrix as the error measure than the mere misclassification error, which corresponds to the probability for some classifier h to err for its prediction on \mathbf{x} . The information as to what is the probability for an instance of class p to be classified into class q (with $p \neq q$) by some predictor is indeed crucial in some applications (think of the difference between false-negative and false-positive predictions in a diagnosis automated system). To the best of our knowledge, we are the first to propose a generalization bound on the confusion matrix in the PAC-Bayesian framework. The result that we propose heavily relies on a matrix concentration inequality for sums of random matrices introduced by Tropp (2011). One may anticipate that generalization bounds for the confusion matrix may also be obtained in other framework than the PAC-Bayesian framework (e.g. uniform stability, online learning).

The rest of this paper is organized as follows. Section 2 introduces the setting of multi-class learning and some of the basic notation used throughout the paper. Section 3 briefly recalls the folk PAC-Bayes bound as introduced in McAllester (2003). In Section 4, we present the main contribution of this paper, our PAC-Bayes bound on the confusion matrix, followed by its proof in Section 5. We discuss some future works in Section 6.

2 Setting and Notations

This section presents the general setting that we consider and the different tools that we will make use of.

2.1 General Problem Setting

We consider classification tasks over the *input space* $X \subseteq \mathbb{R}^d$ of dimension d . The *output space* is denoted by $Y = \{1, \dots, Q\}$, where Q is the number of classes. The learning sample is denoted by $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ where each example is drawn *i.i.d.* from a fixed —but unknown— probability distribution \mathfrak{D} defined over $X \times Y$. \mathfrak{D}_m denotes the distribution of a m -sample. $\mathcal{F} \subseteq \mathbb{R}^X$ is a family of classifiers $f : X \rightarrow Y$. \mathfrak{P} and \mathfrak{Q} are respectively the *prior* and the *posterior* distributions over \mathcal{F} . Given the prior distribution \mathfrak{P} and the training set S , the learning process consists in finding the posterior distribution \mathfrak{Q} leading to a good generalization.

Since we make use of the prior distribution \mathfrak{P} on \mathcal{F} , a PAC-Bayes generalization bound depends on the Kullback-Leibler divergence (KL-divergence):

$$KL(\mathfrak{Q} \parallel \mathfrak{P}) = \mathbb{E}_{f \sim \mathfrak{Q}} \log \frac{\mathfrak{Q}(f)}{\mathfrak{P}(f)}. \quad (1)$$

The function $\text{sign}(x)$ is equal to $+1$ if $x \geq 0$ and -1 otherwise. The indicator function $\mathbb{I}(x)$ is equal to 1 if x is true and 0 otherwise.

2.2 Conventions and Basics on Matrices

Throughout the paper we consider only real-valued square matrices \mathbf{C} of order Q (the number of classes). ${}^t\mathbf{C}$ is the transpose of the matrix \mathbf{C} , \mathbf{Id}_Q denotes the identity matrix of size Q and $\mathbf{0}$ is the zero matrix.

The results given in this paper are based on a concentration inequality of Tropp (2011) for a sum of random self-adjoint matrices. In the case when a matrix is not self-adjoint and is real-valued, we use the dilation of such a matrix, given in Paulsen (2002), which is defined as follows:

$$\mathcal{S}(\mathbf{C}) \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{0} & \mathbf{C} \\ {}^t\mathbf{C} & \mathbf{0} \end{pmatrix}. \quad (2)$$

The symbol $\|\cdot\|$ corresponds to the *operator norm* also called the *spectral norm* since it returns the largest singular value of its argument, which is defined by:

$$\|\mathbf{C}\| = \max\{\lambda_{\max}(\mathbf{C}), -\lambda_{\min}(\mathbf{C})\}, \quad (3)$$

where λ_{\max} and λ_{\min} are respectively the algebraic maximum and minimum singular value of \mathbf{C} . Note that the dilation preserves spectral information, so we have:

$$\lambda_{\max}(\mathcal{S}(\mathbf{C})) = \|\mathcal{S}(\mathbf{C})\| = \|\mathbf{C}\|. \quad (4)$$

An important property of the operator is the following:

$$\forall a \in \mathbb{R}, \|a \cdot \mathbf{C}\| = |a| \cdot \|\mathbf{C}\|. \quad (5)$$

3 The Usual PAC-Bayes Theorem

In this section, we recall the main PAC-Bayesian bound in binary classification case as presented in McAllester (2003); Seeger and Seeger (2002); Langford (2005). The set of labels we consider is $Y = \{-1; 1\}$ (with $Q = 2$) and, for each classifier $f \in \mathcal{F}$, the predicted output of $\mathbf{x} \in X$ is given by $\text{sign}(f(\mathbf{x}))$. The true risk $R(f)$ and the empirical error $R_S(f)$ of f are defined as:

$$R(f) \stackrel{\text{def}}{=} \mathbb{E}_{(\mathbf{x}, y) \sim \mathfrak{D}} \mathbb{I}(f(\mathbf{x}) \neq y) \quad ; \quad R_S(f) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) \neq y_i).$$

The learner's aim is to choose a posterior distribution \mathfrak{Q} on \mathcal{F} such that the risk of the \mathfrak{Q} -weighted majority vote (also called the Bayes classifier) $B_{\mathfrak{Q}}$ is as small as possible. $B_{\mathfrak{Q}}$ is defined by:

$$B_{\mathfrak{Q}}(\mathbf{x}) = \text{sign}[\mathbb{E}_{f \sim \mathfrak{Q}} f(\mathbf{x})].$$

The true risk $R(B_{\mathfrak{Q}})$ and the empirical error $R_S(B_{\mathfrak{Q}})$ of the Bayes classifier are defined as the probability that it commits an error on an example:

$$R(B_{\mathfrak{Q}}) \stackrel{\text{def}}{=} \mathbb{P}_{(\mathbf{x}, y) \sim \mathfrak{D}} (B_{\mathfrak{Q}}(\mathbf{x}) \neq y). \quad (6)$$

However, the PAC-Bayes approach does not directly bound the risk of $B_{\mathfrak{Q}}$. Instead, it bounds the risk of the stochastic Gibbs classifier $G_{\mathfrak{Q}}$ which predicts the label of $\mathbf{x} \in X$ by first drawing f according to \mathfrak{Q} and then returning $f(\mathbf{x})$. The true risk $R(G_{\mathfrak{Q}})$ and the empirical error $R_S(G_{\mathfrak{Q}})$ of $G_{\mathfrak{Q}}$ are therefore:

$$R(G_{\mathfrak{Q}}) = \mathbb{E}_{f \sim \mathfrak{Q}} R(f) \quad ; \quad R_S(G_{\mathfrak{Q}}) = \mathbb{E}_{f \sim \mathfrak{Q}} R_S(f). \quad (7)$$

Note that in this setting, if $B_{\mathfrak{Q}}$ misclassifies \mathbf{x} , then at least half of the classifiers (under \mathfrak{Q}) commit an error on \mathbf{x} . Hence, we directly have: $R(B_{\mathfrak{Q}}) \leq 2R(G_{\mathfrak{Q}})$. Thus, an upper bound on $R(G_{\mathfrak{Q}})$ gives rise to an upper bound on $R(B_{\mathfrak{Q}})$.

We present the PAC-Bayes theorem which gives a bound on the error of the stochastic Gibbs classifier.

Theorem 1 (*i.i.d.* binary classification PAC-Bayes Bound). *For any \mathfrak{D} , any \mathcal{F} , any \mathfrak{P} of support \mathcal{F} , any $\delta \in (0, 1]$, we have,*

$$\mathbb{P}_{S \sim \mathfrak{D}_m} \left(\forall \mathfrak{Q} \text{ on } \mathcal{F}, \text{kl}(R_S(G_{\mathfrak{Q}}), R(G_{\mathfrak{Q}})) \leq \frac{1}{m} \left[KL(\mathfrak{Q} \parallel \mathfrak{P}) + \ln \frac{\xi(m)}{\delta} \right] \right) \geq 1 - \delta,$$

where $\text{kl}(a, b) \stackrel{\text{def}}{=} a \ln \frac{a}{b} + (1 - a) \ln \frac{1-a}{1-b}$, and $\xi \stackrel{\text{def}}{=} \sum_{i=0}^m \binom{m}{i} (i/m)^i (1 - i/m)^{m-i}$.

We now provide a novel PAC-Bayes bound in the context of multi-class classification by considering the confusion matrix as an error measure.

4 Multiclass PAC-Bayes Bound

4.1 Definitions and Setting

As said earlier, we focus on multi-class classification. The output space is $Y = \{1, \dots, Q\}$, with $Q > 2$. We only consider learning algorithms acting on learning sample $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ where each example is drawn *i.i.d* according to \mathfrak{D} , such that $|S| \geq Q$ and $m_{y_j} \geq 1$ for every class $y_j \in Y$, where m_{y_j} is the number of examples of real class y_j . In the context of multi-class classification, an error measure can be the *confusion matrix*. Concretely, for a given classifier $f \in \mathcal{F}$ and a sample $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \sim \mathfrak{D}_m$, the *empirical confusion matrix* $\mathbf{D}_S^f = (\hat{d}_{pq})_{1 \leq p, q \leq Q}$ of f is defined as follows:

$$\forall(p, q), \hat{d}_{pq} \stackrel{\text{def}}{=} \sum_{i=1}^m \frac{1}{m_{y_i}} \mathbb{I}(f(\mathbf{x}_i) = q) \mathbb{I}(y_i = p).$$

The *true confusion matrix* $\mathbf{D}^f = (d_{pq})_{1 \leq p, q \leq Q}$ of f over \mathfrak{D} corresponds to:

$$\begin{aligned} \forall(p, q), d_{pq} &\stackrel{\text{def}}{=} \mathbb{E}_{\mathbf{x}|y=p} \mathbb{I}(f(\mathbf{x}) = q) \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathfrak{D}}(f(\mathbf{x}) = q | p = y). \end{aligned}$$

If f correctly classifies every example of the sample S , then all the elements of the confusion matrix are 0, except for the diagonal ones which correspond to the correctly classified examples. Hence the more there are non-zero elements in a confusion matrix outside the diagonal, the more the classifier is prone to err. Recall that in a learning process the objective is to learn a classifier $f \in \mathcal{F}$ with a low true error (*i.e.* with good generalization guarantees), we are thus only interested in the errors of f . Our objective is then to find f leading to a confusion matrix with the more zero elements outside the diagonal. Therefore, we propose to consider a different kind of confusion matrix by discarding the diagonal values. The only non-zero elements of the new confusion matrix correspond to the examples that are misclassified by f . For all $f \in \mathcal{F}$ we define the empirical and true confusion matrices of f by respectively $\mathbf{C}_S^f = (\hat{c}_{pq})_{1 \leq p, q \leq Q}$ and $\mathbf{C}^f = (c_{pq})_{1 \leq p, q \leq Q}$ such that:

$$\forall(p, q), \hat{c}_{pq} \stackrel{\text{def}}{=} \begin{cases} 0 & \text{if } q = p \\ \hat{d}_{pq} & \text{otherwise,} \end{cases} \quad (8)$$

$$\forall(p, q), c_{pq} \stackrel{\text{def}}{=} \begin{cases} 0 & \text{if } q = p \\ d_{pq} = \mathbb{P}_{(\mathbf{x}, y) \sim \mathfrak{D}}(f(\mathbf{x}) = q | p = y) & \text{otherwise.} \end{cases} \quad (9)$$

Note that if f correctly classifies every example of a given sample S , then the empirical confusion matrix \mathbf{C}_S^f is equal to $\mathbf{0}$. Similarly, if f is a perfect classifier over the distribution \mathfrak{D} , then the true confusion matrix is equal to $\mathbf{0}$. Aiming at controlling the confusion matrix of a classifier is therefore a relevant task. More precisely, one may aim at a confusion matrix that is ‘small’, where ‘small’ means as close to $\mathbf{0}$ as possible. As we shall see, the size of a confusion matrix will be measured by its operator norm.

4.2 Main Result: PAC-Bayes Bound on the Confusion Matrix of the Gibbs Classifier

Our main result is a PAC-Bayes generalization bound over the Gibbs classifier G_Ω in this particular context, where the empirical and true error measures are respectively given by the confusion matrices from (8) and (9). In this case, we can define the true and the empirical confusion matrices of G_Ω respectively by:

$$\mathbf{C}^{G_\Omega} = \mathbb{E}_{f \sim \Omega} \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_S^f ; \quad \mathbf{C}_S^{G_\Omega} = \mathbb{E}_{f \sim \Omega} \mathbf{C}_S^f.$$

Given $f \sim \Omega$ and a sample $S \sim \mathfrak{D}_m$, our objective is to bound the difference between \mathbf{C}^{G_Ω} and $\mathbf{C}_S^{G_\Omega}$, the true and empirical errors of the Gibbs classifier. The structure that we will consider in the space of confusion matrices is the one induced by the operator norm (Equation (3)) on the set of matrices. This norm will allow us to formally relate the true and empirical confusion matrices of the Gibbs classifier and it also will allow us to provide a bound on the size $\|\mathbf{C}^{G_\Omega}\|$ of the true confusion matrix.

Here is our main result.

Theorem 2. Let $X \subseteq \mathbb{R}^d$ be the input space, $Y = \{1, \dots, Q\}$ the output space, \mathfrak{D} a distribution over $X \times Y$ (with \mathfrak{D}_m the distribution of a m -sample) and \mathcal{F} a family of classifiers from X to Y . Then for every prior distribution \mathfrak{P} over \mathcal{F} and any $\delta \in (0, 1]$, we have:

$$\mathbb{P}_{S \sim \mathfrak{D}_m} \left\{ \forall \Omega \text{ on } \mathcal{F}, \|\mathbf{C}_S^{G_\Omega} - \mathbf{C}^{G_\Omega}\| \leq \sqrt{\frac{8Q}{m_- - 8Q} \left[KL(\Omega \| \mathfrak{P}) + \ln \left(\frac{m_-}{4\delta} \right) \right]} \right\} \geq 1 - \delta,$$

where $m_- = \min_{y=1, \dots, Q} m_y$ corresponds to the minimal number of examples from S which belong to the same class.

Proof. Deferred to Section 5. □

Note that, for all $y \in Y$, we need the following hypothesis: $m_y > 8$, which is not too strong a limitation.

Finally, we rewrite Theorem 2 to have the size of the confusion matrix under consideration.

Corollary 1. We consider the hypothesis of the Theorem 2. We have:

$$\mathbb{P}_{S \sim \mathfrak{D}_m} \left\{ \forall \Omega \text{ on } \mathcal{F}, \|\mathbf{C}^{G_\Omega}\| \leq \|\mathbf{C}_S^{G_\Omega}\| + \sqrt{\frac{8Q}{m_- - 8Q} \left[KL(\Omega \| \mathfrak{P}) + \ln \left(\frac{m_-}{4\delta} \right) \right]} \right\} \geq 1 - \delta.$$

Proof. By application of the reverse triangle inequality $\|\mathbf{A}\| - \|\mathbf{B}\| \leq \|\mathbf{A} - \mathbf{B}\|$ to Theorem 2. □

For a fixed prior \mathfrak{P} on \mathcal{F} , both Theorem 2 and Corollary 1 yield a bound on the estimation (through the operator norm) of the true confusion matrix of the Gibbs classifier over all¹ posterior distribution Ω on \mathcal{F} , though this is more explicit in the corollary. Let the number of classes Q be a constant, then the true risk is upper-bounded by the empirical risk of the Gibbs classifier and a term depending on the number of training examples, especially on the value m_- which corresponds to the minimal quantity of examples that belong to the same class. This means that the larger m_- , the closer the empirical confusion matrix of the Gibbs classifier to its true matrix. These bounds use first-order information and vary as $O(1/\sqrt{m_-})$, which is a typical rate of bounds not using second-order information.

5 Proof of Theorem 2

This section gives the formal proof of Theorem 2. We first introduce a concentration inequality for a sum of random square matrices. This allows us to deduce the PAC-Bayes generalization bound for confusion matrices by following the same “three step process” as the one given in McAllester (2003); Seeger and Seeger (2002); Langford (2005) for the classic PAC-Bayesian bound.

5.1 Concentration Inequality for the Confusion Matrix

The main result of our work is based on the following corollary of a result on the concentration inequality for a sum of self-adjoint matrices given by Tropp (2011) (see Theorem 3 in Appendix) – this theorem generalizes Hoeffding’s inequality to the case self-adjoint random matrices. The purpose of the following corollary is to restate the Theorem 3 so that it carries over to matrices that are not self-adjoint. It is central to us to have such a result as the matrices we are dealing with, namely confusion matrices, are rarely symmetric.

Corollary 2. Consider a finite sequence $\{\mathbf{M}_i\}$ of independent, random, square matrices of order Q , and let $\{a_i\}$ be a sequence of fixed scalars. Assume that each random matrix satisfies $\mathbb{E}_i \mathbf{M}_i = \mathbf{0}$ and $\|\mathbf{M}_i\| \leq a_i$ almost surely. Then, for all $\epsilon \geq 0$,

$$\mathbb{P} \left\{ \left\| \sum_i \mathbf{M}_i \right\| \geq \epsilon \right\} \leq 2 \cdot Q \cdot \exp \left(\frac{-\epsilon^2}{8\sigma^2} \right), \quad (10)$$

where $\sigma^2 \stackrel{\text{def}}{=} \sum_i a_i^2$.

¹This includes any Ω chosen by the learner after observing S .

Proof. We want to verify the hypothesis given in Theorem 3 in order to apply it.

Let $\{\mathbf{M}_i\}$ be a finite sequence of independent, random, square matrices of order Q such that $\mathbb{E}_i \mathbf{M}_i = \mathbf{0}$ and let $\{a_i\}$ be a sequence of fixed scalars such that $\|\mathbf{M}_i\| \leq a_i$. We consider the sequence $\{\mathcal{S}(\mathbf{M}_i)\}$ of random self-adjoint matrices with dimension $2Q$. By the definition of the dilation, we directly obtain $\mathbb{E}_i \mathcal{S}(\mathbf{M}_i) = \mathbf{0}$.

From Equation (4), the dilation preserves the spectral information. Thus, on the one hand, we have:

$$\left\| \sum_i \mathbf{M}_i \right\| = \lambda_{\max} \left(\mathcal{S} \left(\sum_i \mathbf{M}_i \right) \right) = \lambda_{\max} \left(\sum_i \mathcal{S}(\mathbf{M}_i) \right).$$

On the other hand, we have:

$$\|\mathbf{M}_i\| = \|\mathcal{S}(\mathbf{M}_i)\| = \lambda_{\max}(\mathcal{S}(\mathbf{M}_i)) \leq a_i.$$

To assure the hypothesis $\mathcal{S}(\mathbf{M}_i)^2 \preceq \mathbf{A}_i^2$, we need to find a suitable sequence of fixed self-adjoint matrices $\{\mathbf{A}_i\}$ of dimension $2Q$ (where \preceq refers to the semidefinite order on self-adjoint matrices). Indeed, it suffices to construct a diagonal matrix defined as $\lambda_{\max}(\mathcal{S}(\mathbf{M}_i)) \mathbf{Id}_{2Q}$ for ensuring $\mathcal{S}(\mathbf{M}_i)^2 \preceq (\lambda_{\max}(\mathcal{S}(\mathbf{M}_i)) \mathbf{Id}_{2Q})^2$. More precisely, since for every i we have $\lambda_{\max}(\mathcal{S}(\mathbf{M}_i)) \leq a_i$, we fix \mathbf{A}_i as a diagonal matrix with a_i on the diagonal, *i.e.* $\mathbf{A}_i \stackrel{\text{def}}{=} a_i \mathbf{Id}_{2Q}$, with $\|\sum_i \mathbf{A}_i^2\| = \sum_i a_i^2 = \sigma^2$.

Finally, we can invoke Theorem 3 to obtain the concentration inequality (10). \square

In order to make use of this corollary, we rewrite confusion matrices as sums of example-based confusion matrices. That is, for each example $(\mathbf{x}_i, y_i) \in S$, we define its empirical confusion matrix by $\mathbf{C}_i^f = (\hat{c}_{pq}(i))_{1 \leq p, q \leq Q}$ as follows:

$$\forall p, q, \hat{c}_{pq}(i) \stackrel{\text{def}}{=} \begin{cases} 0 & \text{if } q = p \\ \frac{1}{m_{y_i}} \mathbb{I}(f(\mathbf{x}) = q) \mathbb{I}(p = y_i) & \text{otherwise.} \end{cases}$$

where m_{y_i} is the number of examples of class $y_i \in Y$ belonging to S . Given an example $(\mathbf{x}_i, y_i) \in S$, the example-based confusion matrix contains at most one non zero-element when f misclassifies (\mathbf{x}_i, y_i) . In the same way, when f correctly classifies (\mathbf{x}_i, y_i) then the example-based confusion matrix is equal to $\mathbf{0}$. Concretely, for every sample $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ and every $f \in \mathcal{F}$, our error measure is then $\mathbf{C}_S^f = \sum_{i=1}^m \mathbf{C}_i^f$. It naturally appears that we penalize only when f errs.

We further introduce the random square matrices $\mathbf{C}_i'^f$:

$$\mathbf{C}_i'^f = \mathbf{C}_i^f - \mathbb{E}_{S \sim \mathcal{D}_m} \mathbf{C}_i^f, \quad (11)$$

which verify $\mathbb{E}_i \mathbf{C}_i'^f = \mathbf{0}$.

We have yet to find a suitable a_i for a given $\mathbf{C}_i'^f$. Let λ_{\max_i} be the maximum singular value of $\mathbf{C}_i'^f$. It is easy to verified that $\lambda_{\max_i} \leq \frac{1}{m_{y_i}}$. Thus, for all i we fix a_i equal to $\frac{1}{m_{y_i}}$.

Finally, with the introduced notations, Corollary 2 leads to the following concentration inequality:

$$\mathbb{P} \left\{ \left\| \sum_{i=1}^m \mathbf{C}_i'^f \right\| \geq \epsilon \right\} \leq 2.Q \cdot \exp \left(\frac{-\epsilon^2}{8\sigma^2} \right). \quad (12)$$

This inequality (12) allows us to demonstrate our Theorem 2 by following the process of McAllester (2003); Seeger and Seeger (2002); Langford (2005).

5.2 “Three Step Proof” Of Our Bound

First, thanks to concentration inequality (12), we prove the following lemma.

Lemma 1. *Let Q be the size of \mathbf{C}_S^f and $\mathbf{C}_i'^f = \mathbf{C}_i^f - \mathbb{E}_{S \sim \mathcal{D}_m} \mathbf{C}_i^f$ defined as in (11). Then the following bound holds for any $\delta \in (0, 1]$:*

$$\mathbb{P}_{S \sim \mathcal{D}_m} \left\{ \mathbb{E}_{f \sim \mathfrak{P}} \left[\exp \left(\frac{1 - 8\sigma^2}{8\sigma^2} \left\| \sum_{i=1}^m \mathbf{C}_i'^f \right\|^2 \right) \right] \leq \frac{2Q}{8\sigma^2\delta} \right\} \geq 1 - \delta$$

Proof. For readability reasons, we note $\mathbf{C}'_S^f = \sum_{i=1}^m \mathbf{C}'_i^f$. If Z is a real valued random variable so that $\mathbb{P}(Z \geq z) \leq k \exp(-n \cdot g(z))$ with $g(z)$ non-negative, non-decreasing and k a constant, then $\mathbb{P}(\exp((n-1)g(Z)) \geq \nu) \leq \min(1, k\nu^{-n/(n-1)})$. We apply this to the concentration inequality (12). Choosing $g(z) = z^2$ (non-negative), $z = \epsilon$, $n = \frac{1}{8\sigma^2}$ and $k = 2Q$, we obtain the following result:

$$\mathbb{P}\left\{\exp\left(\frac{1-8\sigma^2}{8\sigma^2}\|\mathbf{C}'_S^f\|\right) \geq \nu\right\} \leq \min(1, 2Q\nu^{-1/(1-8\sigma^2)}).$$

Note that $\exp\left(\frac{1-8\sigma^2}{8\sigma^2}\|\mathbf{C}'_S^f\|\right)$ is always non-negative. Hence it allows us to compute its expectation as:

$$\begin{aligned} \mathbb{E}\left[\exp\left(\frac{1-8\sigma^2}{8\sigma^2}\|\mathbf{C}'_S^f\|\right)\right] &= \int_0^\infty \mathbb{P}\left\{\exp\left(\frac{1-8\sigma^2}{8\sigma^2}\|\mathbf{C}'_S^f\|\right) \geq \nu\right\} d\nu \\ &\leq 2Q + \int_1^\infty 2Q\nu^{-1/(1-8\sigma^2)} d\nu \\ &= 2Q - 2Q \frac{1-8\sigma^2}{8\sigma^2} \left[\nu^{-8\sigma^2/(1-8\sigma^2)}\right]_1^\infty \\ &= 2Q + 2Q \frac{1-8\sigma^2}{8\sigma^2} \\ &= \frac{2Q}{8\sigma^2}. \end{aligned}$$

For a given classifier $f \in \mathcal{F}$, we have:

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[\exp\left(\frac{1-8\sigma^2}{8\sigma^2}\|\mathbf{C}'_S^f\|\right) \right] \leq \frac{2Q}{8\sigma^2}. \quad (13)$$

Then, if \mathfrak{P} is a probability distribution over \mathcal{F} , Equation (13) implies that:

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[\mathbb{E}_{f \sim \mathfrak{P}} \exp\left(\frac{1-8\sigma^2}{8\sigma^2}\|\mathbf{C}'_S^f\|\right) \right] \leq \frac{2Q}{8\sigma^2}. \quad (14)$$

Using Markov's inequality², we obtain the result of the lemma. \square

The second step to prove Theorem 2 is to use the shift given in McAllester (2003). We recall this result in the following lemma.

Lemma 2 (McAllester (2003)). *Given the Kullback-Leibler divergence³ $KL(\mathfrak{Q} \parallel \mathfrak{P})$ between two distributions \mathfrak{P} and \mathfrak{Q} and let $g(\cdot)$ be a function, we have:*

$$\mathbb{E}_{a \sim \mathfrak{Q}} [g(b)] \leq KL(\mathfrak{Q} \parallel \mathfrak{P}) + \ln \mathbb{E}_{x \sim \mathfrak{P}} [\exp(g(b))].$$

Proof. See McAllester (2003). \square

Recall that $\mathbf{C}'_S^f = \sum_{i=1}^m \mathbf{C}'_i^f$. With $g(b) = \frac{1-8\sigma^2}{8\sigma^2}b^2$ and $b = \|\mathbf{C}'_S^f\|$, Lemma 2 implies:

$$\mathbb{E}_{f \sim \mathfrak{Q}} \left[\frac{1-8\sigma^2}{8\sigma^2} \|\mathbf{C}'_S^f\|^2 \right] \leq KL(\mathfrak{Q} \parallel \mathfrak{P}) + \ln \mathbb{E}_{f \sim \mathfrak{P}} \left[\exp\left(\frac{1-8\sigma^2}{8\sigma^2} \|\mathbf{C}'_S^f\|^2\right) \right]. \quad (15)$$

The last step that completes the proof of Theorem 2 consists in applying the result we obtained in Lemma 1 to Equation (15). Then, we have:

$$\mathbb{E}_{f \sim \mathfrak{Q}} \left[\frac{1-8\sigma^2}{8\sigma^2} \|\mathbf{C}'_S^f\|^2 \right] \leq KL(\mathfrak{Q} \parallel \mathfrak{P}) + \ln \frac{2Q}{8\sigma^2\delta}. \quad (16)$$

Since $g(\cdot)$ is clearly convex, we apply Jensen's inequality⁴ to (16). Then, with probability at least $1 - \delta$ over S , and for every distribution \mathfrak{Q} on \mathcal{F} , we have:

$$\left(\mathbb{E}_{f \sim \mathfrak{Q}} \|\mathbf{C}'_S^f\| \right)^2 \leq \frac{8\sigma^2}{1-8\sigma^2} \left(KL(\mathfrak{Q} \parallel \mathfrak{P}) + \ln \frac{2Q}{8\sigma^2\delta} \right). \quad (17)$$

Since $\mathbf{C}'_S^f = \sum_{i=1}^m [\mathbf{C}_i^f - \mathbb{E}_{S \sim \mathcal{D}_m} \mathbf{C}_i^f]$, then the bound (17) is quite similar to the one given in Theorem 2.

We present in the next section, the calculations leading to our PAC-Bayesian generalization bound.

²see Theorem 4 in Appendix.

³The KL-divergence is defined in Equation (1).

⁴see Theorem 5 in Appendix.

5.3 Simplification

We first compute the variance parameter $\sigma^2 = \sum_{i=1}^m a_i^2$. For that purpose, in Section 5.1 we showed that for each $i \in \{1, \dots, m\}$, we can choose $a_i = \frac{1}{m_{y_i}}$, where y_i is the class of the i -th example and m_{y_i} is the number of examples of class y_i . Thus we have:

$$\sigma^2 = \sum_{i=1}^m \frac{1}{m_{y_i}^2} = \sum_{y=1}^Q \sum_{i: y_i=y} \frac{1}{m_y^2} = \sum_{y=1}^Q \frac{1}{m_y}.$$

For sake of simplification of Equation (17) and since the term on the right side of this equation is an increasing function with respect to σ^2 , we propose to upper-bound σ^2 :

$$\sigma^2 = \sum_{y=1}^Q \frac{1}{m_y} \leq \frac{Q}{\min_{y=1, \dots, Q} m_y}. \quad (18)$$

Let $m_- \stackrel{\text{def}}{=} \min_{y=1, \dots, Q} m_y$, then using Equation (18), we obtain the following bound from Equation (17):

$$\left(\mathbb{E}_{f \sim \Omega} [\| \mathbf{C}'_S^f \|] \right)^2 \leq \frac{8Q}{m_- - 8Q} \left(KL(\Omega \| \mathfrak{P}) + \ln \frac{m_-}{4\delta} \right).$$

Then:

$$\mathbb{E}_{f \sim \Omega} [\| \mathbf{C}'_S^f \|] \leq \sqrt{\frac{8Q}{m_- - 8Q} \left(KL(\Omega \| \mathfrak{P}) + \ln \frac{m_-}{4\delta} \right)}. \quad (19)$$

It remains to replace $\mathbf{C}'_S^f = \sum_{i=1}^m [\mathbf{C}_i^f - \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_i^f]$. Recall that $\mathbf{C}^{G_\Omega} = \mathbb{E}_{f \sim \Omega} \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_S^f$ and $\mathbf{C}_S^{G_\Omega} = \mathbb{E}_{f \sim \Omega} \mathbf{C}_S^f$, we obtain:

$$\begin{aligned} \mathbb{E}_{f \sim \Omega} [\| \mathbf{C}'_S^f \|] &= \mathbb{E}_{f \sim \Omega} \left[\left\| \sum_{i=1}^m [\mathbf{C}_i^f - \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_i^f] \right\| \right] \\ &= \mathbb{E}_{f \sim \Omega} \left[\left\| \sum_{i=1}^m [\mathbf{C}_i^f] - \sum_{i=1}^m [\mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_i^f] \right\| \right] \\ &= \mathbb{E}_{f \sim \Omega} \left[\left\| \mathbf{C}_S^f - \mathbb{E}_{S \sim \mathfrak{D}_m} \left[\sum_{i=1}^m \mathbf{C}_i^f \right] \right\| \right] \\ &= \mathbb{E}_{f \sim \Omega} [\| \mathbf{C}_S^f - \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_S^f \|] \\ &\geq \mathbb{E}_{f \sim \Omega} [\| \mathbf{C}_S^f - \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_S^f \|] \\ &= \mathbb{E}_{f \sim \Omega} \mathbf{C}_S^f - \mathbb{E}_{f \sim \Omega} \mathbb{E}_{S \sim \mathfrak{D}_m} \mathbf{C}_S^f \\ &= \| \mathbf{C}_S^{G_\Omega} - \mathbf{C}^{G_\Omega} \|. \end{aligned} \quad (20)$$

By substituting the left part of the inequality (19) with the term (20), we find the bound of our Theorem 2.

6 Discussion and Future Work

This work gives rise to many interesting questions, among which the following ones.

In the case of the classical binary PAC-Bayes framework, it is easy to show that the true error of the Bayes classifier (6) and the one of the Gibbs classifier (7) are related by the following inequality:

$$R(B_\Omega) \leq 2R(G_\Omega).$$

One may notice that we do not immediately have a similar result for our confusion matrix setting. This question is out of the scope of the present paper and the proof of such a relation between the confusion

matrix-based errors of the Bayes classifier and of the Gibbs classifier for this framework is left for future work.

Other perspectives will be focused on instantiating our bound given in Theorem 2 for specific multi-class frameworks, such as multi-class SVM Weston and Watkins (1998); Crammer and Singer (2002); Lee et al. (2004) and multi-class boosting (AdaBoost.MH/AdaBoost.MR Schapire and Singer (2000), SAMME Zhu et al. (2009), AdaBoost.MM Mukherjee and Schapire (2011)). Taking advantage of our theorem while using the confusion matrices, may allow us to derive new generalization bounds for these methods.

Additionally, we are interested in seeing how effective learning methods may be derived from the risk bound we propose. For instance, in the binary PAC-Bayes setting, the algorithm MinCq proposed by Laviolette et al. (2011) minimizes a bound depending on the first two moments of the margin of the Q -weighted majority vote. From our Theorem 2 and with a similar study, we would like to design a new multi-class learning algorithm and observe how sound such an algorithm could be. This would probably require the derivation of a Cantelli-Tchebycheff deviation inequality in the matrix case.

Besides, it might be very interesting to see how the noncommutative/matrix concentration inequalities provided by Tropp (2011) might be of some use for other kinds of learning problem such as multi-label classification, label ranking problems or structured prediction issues.

Finally, the question of extending the present work to the analysis of algorithms learning (possibly infinite-dimensional) operators as Abernethy et al. (2009) is also very exciting.

7 Conclusion

In this paper, we propose a new PAC-Bayesian generalization bound that applies in the multi-class classification setting. The originality of our contribution is that we consider the confusion matrix as an error measure. Coupled with the use of the operator norm on matrices, we are capable of providing generalization bound on the ‘size’ of confusion matrix (with the idea that the smaller the norm of the confusion matrix of the learned classifier, the better it is for the classification task at hand). The derivation of our result takes advantage of the concentration inequality proposed by Tropp (2011) for the sum of random self-adjoint matrices, that we directly adapt to square matrices which are not self-adjoint.

The main results are presented in Theorem 2 and Corollary 1. The bound in Theorem 2 is given on the difference between the true risk of the Gibbs classifier and its empirical error. While the one given in Corollary 1 upper-bounds the risk of the Gibbs classifier by its empirical error.

An interesting point is that our bound depends on the minimal quantity m_- of training examples belonging to the same class, for a given number of classes. If this value increases, *i.e.* if we have a lot of training examples, then the empirical confusion matrix of the Gibbs classifier tends to be close to its true confusion matrix. A point worth noting is that the bound varies as $O(1/\sqrt{m_-})$, which is a typical rate in bounds not using second-order information.

The present work gives rise to a few algorithmic and theoretical questions that we have discussed in the previous section.

Appendix

Theorem 3 (Concentration Inequality for Random Matrices Tropp (2011)). *Consider a finite sequence $\{\mathbf{M}_i\}$ of independent, random, self-adjoint matrices with dimension Q , and let $\{\mathbf{A}_i\}$ be a sequence of fixed self-adjoint matrices. Assume that each random matrix satisfies $\mathbb{E}\mathbf{M}_i = \mathbf{0}$ and $\mathbf{M}_i^2 \preceq \mathbf{A}_i^2$ almost surely. Then, for all $\epsilon \geq 0$,*

$$\mathbb{P} \left\{ \lambda_{\max} \left(\sum_i \mathbf{M}_i \right) \geq \epsilon \right\} \leq Q \cdot \exp \left(\frac{-\epsilon^2}{8\sigma^2} \right),$$

where $\sigma^2 \stackrel{\text{def}}{=} \left\| \sum_i \mathbf{A}_i^2 \right\|$ and \preceq refers to the semidefinite order on self-adjoint matrices.

Theorem 4 (Markov’s inequality). *Let Z be a random variable and $z \geq 0$, then:*

$$\mathbb{P}(|Z| \geq z) \leq \frac{\mathbb{E}(|Z|)}{z}.$$

Theorem 5 (Jensen’s inequality). *Let X be an integrable real-valued random variable and $g(\cdot)$ be a convex function, then:*

$$f(\mathbb{E}[Z]) \leq \mathbb{E}[g(Z)].$$

References

- Abernethy, J., Bach, F., Evgeniou, T., and Vert, J.-P. (2009). A new approach to collaborative filtering: Operator estimation with spectral regularization. *Journal of Machine Learning Research*, 10:803–826.
- Catoni, O. (2004). 4. Gibbs estimators. In *Statistical Learning Theory and Stochastic Optimization*, volume 1851, pages 111–135. Springer.
- Catoni, O. (2007). PAC-bayesian supervised classification: The thermodynamics of statistical learning. *ArXiv e-prints*.
- Crammer, K. and Singer, Y. (2002). On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2:265–292.
- Freund, Y. and Schapire, R. E. (1996). Experiments with a new boosting algorithm. In *In Proceedings of the International Conference on Machine Learning*, pages 148–156.
- Lacasse, A., Laviolette, F., Marchand, M., Germain, P., and Usunier, N. (2007). PAC-bayes bounds for the risk of the majority vote and the variance of the Gibbs classifier. In *Proceedings of the conference on Neural Information Processing Systems (NIPS)*.
- Langford, J. (2005). Tutorial on practical prediction theory for classification. *Journal of Machine Learning Research*, 6:273–306.
- Langford, J., Seeger, M., and Megiddo, N. (2001). An improved predictive accuracy bound for averaging classifiers. In *Proceedings of the International Conference on Machine Learning*, pages 290–297.
- Langford, J. and Shawe-taylor, J. (2002). PAC-bayes & margins. In *Advances in Neural Information Processing Systems 15*, pages 439–446. MIT Press.
- Laviolette, F., Marchand, M., and Roy, J.-F. (2011). From PAC-Bayes Bounds to Quadratic Programs for Majority Votes. In *Proceedings of the International Conference on Machine Learning*.
- Lee, Y., Lin, Y., and Wahba, G. (2004). Multicategory support vector machines, theory, and application to the classification of microarray data and satellite radiance data. *Journal of the American Statistical Association*, 99:67–81.
- McAllester, D. A. (1999a). PAC-bayesian model averaging. In *Proceedings of the annual conference on Computational learning theory (COLT)*, pages 164–170.
- McAllester, D. A. (1999b). Some PAC-bayesian theorems. *Machine Learning*, 37:355–363.
- McAllester, D. A. (2003). Simplified PAC-bayesian margin bounds. In *Proceedings of the annual conference on Computational learning theory (COLT)*, pages 203–215.
- Mukherjee, I. and Schapire, R. E. (2011). A theory of multiclass boosting. *CoRR*, abs/1108.2989.
- Paulsen, V. (2002). *Completely bounded maps and operator algebras*. Cambridge studies in advanced mathematics. Cambridge University Press.
- Schapire, R. E. and Singer, Y. (1999). Improved boosting algorithms using confidence-rated predictions. In *Machine Learning*, pages 80–91.
- Schapire, R. E. and Singer, Y. (2000). BoosTexter: A boosting-based system for text categorization. *Machine Learning*, 39(2/3):135–168.
- Seeger, M. and Seeger, M. (2002). PAC-bayesian generalization error bounds for gaussian process classification. *Journal of Machine Learning Research*, 3:233–269.
- Tropp, J. A. (2011). User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, pages 1–46.
- Weston, J. and Watkins, C. (1998). Multi-class support vector machines.
- Zhu, J., Zou, H., Rosset, S., and Hastie, T. (2009). Multi-class adaboost.